

UNCLASSIFIED

AD 273 207

*Reproduced
by the*

**ARMED SERVICES TECHNICAL INFORMATION AGENCY
ARLINGTON HALL STATION
ARLINGTON 12, VIRGINIA**



UNCLASSIFIED

NOTICE: When government or other drawings, specifications or other data are used for any purpose other than in connection with a definitely related government procurement operation, the U. S. Government thereby incurs no responsibility, nor any obligation whatsoever; and the fact that the Government may have formulated, furnished, or in any way supplied the said drawings, specifications, or other data is not to be regarded by implication or otherwise as in any manner licensing the holder or any other person or corporation, or conveying any rights or permission to manufacture, use or sell any patented invention that may in any way be related thereto.

D1-82-0152

273207

ASTIA

CATALOGED BY

AS AD No. _____

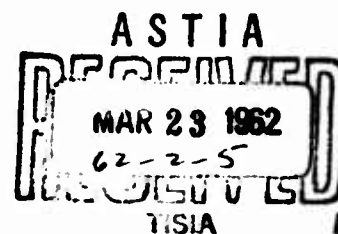
BOEING SCIENTIFIC RESEARCH LABORATORIES

Uniform Distributions Over a Simplex

G. Marsaglia

Mathematics Research

December, 1961



UNIFORM DISTRIBUTIONS OVER A SIMPLEX

by

G. Marsaglia

Mathematical Note No. 250

Mathematics Research Laboratory

BOEING SCIENTIFIC RESEARCH LABORATORIES

December 1961

UNIFORM DISTRIBUTIONS OVER A SIMPLEX

1. Introduction

We are interested in a random assignment of $n + 1$ probabilities p_1, \dots, p_{n+1} , with $p_1 + \dots + p_{n+1} = 1$. By that we mean that the random vector $\pi = (p_1, \dots, p_{n+1})$ is uniformly distributed over the simplex $S_{n+1} = \{(x_1, \dots, x_{n+1}) : x_i \geq 0, \sum_{i=1}^{n+1} x_i = 1\}$, and this, in turn, has a meaning we will make precise in section 2. We will find the joint distribution of the p 's and consider various ways of representing the p 's in terms of familiar random variables. It turns out that the joint distribution of the p 's is very simple when given in terms of $P[p_1 > a_1, \dots, p_{n+1} > a_{n+1}]$ rather than the conventional $P[p_1 < a_1, \dots, p_{n+1} < a_{n+1}]$, in fact, $P[p_1 > a_1, \dots, p_{n+1} > a_{n+1}] = (1 - a_1 - \dots - a_{n+1})^n$.

The principal purpose of the development is to get Theorem 4, which provides a method for generating exponential (and hence normal, using polar coordinates) random variables in a computer -- π is chosen from S_{n+1} , then a random variable z is chosen so that $zp_1, zp_2, \dots, zp_{n+1}$ are independent and exponentially distributed. We will also point out how π may be used to produce points uniformly over the simplex supported by an arbitrary set of $n + 1$ linearly independent points in space. Thus the main purpose of the development is to provide some theory on which certain "Monte Carlo" techniques may be based, but the development will also provide an approach to problems involving uniform order statistics which is different from that usually taken.

2. Uniform Distributions

Let R_n be Euclidean n -space, and let λ be Lebesgue measure. We say that z_1, \dots, z_n are jointly uniformly distributed, or that the vector $\zeta = (z_1, \dots, z_n)$ is uniformly distributed, over a subset $S \subset R_n$ with positive Lebesgue measure, if the probability measure for ζ is proportional to Lebesgue measure -- for λ -measurable B , $P[(z_1, \dots, z_n) \in B] = \lambda(B)/\lambda(S)$. Then z_1, \dots, z_n have a joint density function; it is constant, $1/\lambda(S)$, over S and zero elsewhere. If A is the $n \times n$ non-singular matrix of a linear transformation T , then $T(\zeta)$, $(= \zeta A)$, is uniformly distributed over the image set $T(S)$.

Let G be the simplex supported by k linearly independent points a_1, \dots, a_k in R_k : $G = \{\beta: \beta = \sum_{i=1}^k c_i a_i, c_i \geq 0, \sum_{i=1}^k c_i = 1\}$. We say that $\eta = (y_1, \dots, y_k)$ is uniformly distributed over the simplex G if for some i the vector $(y_1, y_2, \dots, y_{i-1}, y_{i+1}, \dots, y_k)$ is uniformly distributed, in the above sense, over the projection of G onto R_{k-1} formed by suppressing the i th coordinate, i.e., $T_i(G)$, where T_i is the transformation $(y_1, \dots, y_k) \rightarrow (y_1, y_2, \dots, y_{i-1}, y_{i+1}, \dots, y_k)$ taking R_k onto R_{k-1} . This is equivalent to saying that the probability that η fall in a subset H of G is proportional to the surface area of H . For each i , $T_i(\eta)$ will be uniformly distributed in the first sense over $T_i(G)$, except in the nuisance case when one of the coordinates, say x_j , of each point $\xi = (x_1, \dots, x_k)$ of G is constant; $x_j = c$. Then only $T_j(\eta)$ will be uniformly distributed.

We note in particular that if $\pi = (p_1, \dots, p_{n+1})$ is uniformly distributed over

$$S_{n+1} = \{(x_1, \dots, x_{n+1}) : x_i \geq 0, \sum_{i=1}^{n+1} x_i = 1\},$$

then each set of n of the p 's is uniformly distributed over the regular polytope

$$S_n^* = \{(x_1, \dots, x_n) : x_i \geq 0, \sum_{i=1}^n x_i \leq 1\}.$$

If A is the $k \times k$ non-singular matrix of a linear transformation T , then $T(\eta)$, $(=\eta A)$, is uniformly distributed over the simplex $T(G)$. This follows from noting that T is the product of elementary linear transformations, and if y_1, \dots, y_k are uniformly distributed over a simplex, then ay_1, y_2, \dots, y_k ($a \neq 0$), or $y_2, y_1, y_3, \dots, y_k$ or $y_1 + cy_2, y_2, \dots, y_k$ are uniformly distributed over the appropriately transformed simplex. We have in particular:

Theorem 1. If $\pi = (p_1, \dots, p_{n+1})$ is uniformly distributed over the simplex $S_{n+1} = \{(x_1, \dots, x_{n+1}) : x_i \geq 0, \sum_{i=1}^{n+1} x_i = 1\}$ and if a_1, \dots, a_{n+1} are $n+1$ linearly independent points in R_{n+1} , then the random vector $p_1 a_1 + p_2 a_2 + \dots + p_{n+1} a_{n+1}$ is uniformly distributed over G , the simplex supported by the a 's: $G = \{\gamma : \gamma = \sum_{i=1}^{n+1} c_i a_i, c_i \geq 0, \sum_{i=1}^{n+1} c_i = 1\}$.

This theorem provides a practical method for producing random points uniformly over regions with linear boundaries, for example, in certain Monte Carlo procedures, and it also shows how the distributions of the coordinates of a random uniform point from a simplex may be found in terms of the distribution of a linear combination of p_1, \dots, p_{n+1} .

If $\gamma_1, \dots, \gamma_n$ are n linearly independent points in R_n and (p_1, \dots, p_n) is uniformly distributed over $S_n^* = \{(x_1, \dots, x_n): x_i \geq 0, \sum x_i \leq 1\}$, then $p_1\gamma_1 + \dots + p_n\gamma_n$ is uniformly distributed over the set $H = \{\beta: \beta = \sum c_i\gamma_i, c_i \geq 0, \sum c_i \leq 1\}$. Let γ_0 be the origin. Then H may be viewed as the convex polytope supported by $\gamma_0, \gamma_1, \dots, \gamma_n$, i.e.,

$$H = \{\delta: \delta = c_1\gamma_1 + \dots + c_n\gamma_n + c_{n+1}\gamma_0, c_i \geq 0, \sum_{i=1}^{n+1} c_i = 1\}.$$

Hence, if $\beta_1 - \beta_0, \dots, \beta_n - \beta_0$ are linearly independent, and if $\pi = (p_1, \dots, p_n, p_{n+1})$ is uniformly distributed over

$$S_{n+1} = \{(x_1, \dots, x_{n+1}): x_i \geq 0, \sum_{i=1}^{n+1} x_i = 1\},$$

then $p_1\beta_1 + \dots + p_n\beta_n + p_{n+1}\beta_0 = p_1(\beta_1 - \beta_0) + p_2(\beta_2 - \beta_0) + \dots + p_n(\beta_n - \beta_0) + \beta_0$ is uniformly distributed over the convex hull of $\beta_0, \beta_1, \dots, \beta_n$, that set being the translate of the convex hull of $\beta_0 - \beta_0, \beta_1 - \beta_0, \dots, \beta_n - \beta_0$.

Theorem 2. If $\pi = (p_1, \dots, p_{n+1})$ is uniformly distributed over the simplex $S_{n+1} = \{(x_1, \dots, x_{n+1}): x_i \geq 0, \sum x_i = 1\}$ and if $\beta_0, \beta_1, \dots, \beta_n$ are $n+1$ points in R_n such that $\beta_1 - \beta_0, \dots, \beta_n - \beta_0$ are linearly independent, then the random vector $p_1\beta_1 + \dots + p_n\beta_n + p_{n+1}\beta_0$ is uniformly distributed over H , the convex polytope supported by the β 's:

$$H = \{\delta: \delta = d_1\beta_1 + \dots + d_n\beta_n + d_{n+1}\beta_0, d_i \geq 0, \sum_{i=1}^{n+1} d_i = 1\}.$$

The most common application of this result is in the plane -- to choose a point β from the triangle formed by $\beta_0, \beta_1, \beta_2$ we choose (p_1, p_2, p_3) uniformly from S_3 and put $\beta = p_1\beta_1 + p_2\beta_2 + p_3\beta_0$. To choose a point uniformly from a polygonal region in the plane, we divide

the region into triangles, choose a triangle with probability proportional to its area, then choose a point uniformly from that triangle.

If (p_1, \dots, p_n) is uniformly distributed over the polytope $S_n^* = \{(x_1, \dots, x_n): x_i \geq 0, \sum x_i \leq 1\}$, then the vector $\omega = (p_1, p_1 + p_2, \dots, p_1 + p_2 + \dots + p_n)$ is uniformly distributed over $Q_n = \{(y_1, \dots, y_n): 0 \leq y_1 \leq y_2 \leq \dots \leq y_n \leq 1\}$ and hence the components of ω are distributed as the order statistics of independent uniform $[0,1]$ random variables u_1, \dots, u_n . Going in the other direction, if $u_{(1)} \leq u_{(2)} \leq \dots \leq u_{(n)}$ are the order statistics of n independent uniform $[0,1]$ random variables, then the vector $(u_{(1)}, u_{(2)} - u_{(1)}, \dots, u_{(n)} - u_{(n-1)})$ is uniformly distributed over S_n^* and the vector $(u_{(1)}, u_{(2)} - u_{(1)}, \dots, u_{(n)} - u_{(n-1)}, 1 - u_{(n)})$ is uniformly distributed over

$$S_{n+1} = \{(x_1, \dots, x_{n+1}): x_i \geq 0, \sum_{i=1}^{n+1} x_i = 1\}.$$

This relationship provides a practical method for producing π uniformly on S_{n+1} in a computer -- we order a set of n independent uniform random variables and take successive differences.

We will need the following fact: if u_1, \dots, u_n are independent random variables, each uniform on $[0,1]$, and $0 \leq c \leq 1$, then

$$(1) \quad P[u_1 + \dots + u_n < c] = c^n/n!.$$

This is a special case of the well-known result on the distribution of the sum of uniform random variables, but it can be proved by induction

quite easily, since

$$P[u_1 + \dots + u_n < c] = \int_0^1 P[u_1 + \dots + u_n < c | u_n = x] dx = \int_0^c P[u_1 + \dots + u_{n-1} < 1 - x] dx.$$

We note that (1) gives the Lebesgue measure of the region

$$(2) \quad A = \{(x_1, \dots, x_n) : x_1 > 0, \dots, x_n > 0, x_1 + \dots + x_n < c\}.$$

3. The Joint Distribution of p_1, \dots, p_{n+1}

Theorem 3. If $\pi = (p_1, \dots, p_n, p_{n+1})$ is uniformly distributed over $S_{n+1} = \{(x_1, \dots, x_{n+1}) : x_i > 0, \sum_{i=1}^{n+1} x_i = 1\}$, then the joint distribution of the p 's is, for non-negative a_1, \dots, a_{n+1} :

$$P[p_1 > a_1, p_2 > a_2, \dots, p_n > a_n, p_{n+1} > a_{n+1}] = (1 - a_1 - a_2 - \dots - a_{n+1})^n$$

for $a_1 + \dots + a_{n+1} \leq 1$ and zero for $a_1 + \dots + a_{n+1} > 1$.

The proof follows readily from the fact that the two sets

$$A = \{(x_1, \dots, x_n) : x_1 > 0, \dots, x_n > 0, x_1 + \dots + x_n < c\}$$

$$B = \{(y_1, \dots, y_n) : y_1 - a_1 > 0, \dots, y_n - a_n > 0, (y_1 - a_1) + (y_2 - a_2) + \dots + (y_n - a_n) < c\}$$

are congruent, and, by (2) and the remark preceding (2), their common Lebesgue measure is $c^n/n!$, for $0 \leq c \leq 1$. Then

$$\begin{aligned} P[p_1 > a_1, \dots, p_{n+1} > a_{n+1}] &= P[p_1 > a_1, \dots, p_n > a_n, 1 - p_1 - \dots - p_n > a_{n+1}] \\ &= P[p_1 - a_1 > 0, \dots, p_n - a_n > 0, (p_1 - a_1) + \dots + (p_n - a_n) < 1 - a_1 - \dots - a_{n+1}]. \end{aligned}$$

The latter is the probability that (p_1, \dots, p_n) will fall in B , with $c = 1 - a_1 - \dots - a_n - a_{n+1}$, and hence equals $\lambda(B)/\lambda(S_n^*) = c^n$.

Putting the appropriate a 's equal to zero, we have

Corollary. If $\pi = (p_1, \dots, p_n, p_{n+1})$ is uniformly distributed over S_{n+1} , then the joint distribution of $p_1, \dots, p_k, (k \leq n)$, is given by

$$P[p_1 > a_1, \dots, p_k > a_k] = (1 - a_1 - a_2 - \dots - a_k)^n, \quad a_i > 0, a_1 + \dots + a_k \leq 1,$$

and, taking partial derivatives, the joint density of p_1, \dots, p_k is

$$n(n+1) \dots (n-k+1)(1-x_1-\dots-x_k)^{n-k} \quad \text{for } x_i > 0, x_1 + \dots + x_k < 1$$

and zero elsewhere.

We now get this interesting result:

Theorem 4. If $\pi = (p_1, \dots, p_{n+1})$ is uniformly distributed over $S_{n+1} = \{(x_1, \dots, x_{n+1}) : x_i \geq 0, \sum_{i=1}^{n+1} x_i = 1\}$, and if z has density function $x^n e^{-x}/n!$, $x > 0$, and is independent of π , then the random variables

$$zp_1, zp_2, \dots, zp_{n+1}$$

are independent and exponentially distributed, (density e^{-x} , $x \geq 0$).

The proof comes from integrating the conditional probability for given z :

$$\begin{aligned} P[zp_1 > b_1, \dots, zp_{n+1} > b_{n+1}] &= \int_0^\infty P[zp_1 > b_1, \dots, zp_{n+1} > b_{n+1} | z = x] \frac{x^n e^{-x}}{n!} dx \\ &= \int_0^\infty P[p_1 > \frac{b_1}{x}, \dots, p_{n+1} > \frac{b_{n+1}}{x}] \frac{x^n e^{-x}}{n!} dx = \int_{b_1 + \dots + b_{n+1}}^\infty \left(1 - \frac{b_1}{x} - \dots - \frac{b_{n+1}}{x}\right)^n \frac{x^n e^{-x}}{n!} dx. \end{aligned}$$

Letting $c = b_1 + \dots + b_{n+1}$, and $y = x - c$, we have

$$\int_c^\infty \frac{(x-c)^n e^{-x}}{n!} dx = e^{-c} \int_0^\infty \frac{y^n e^{-y}}{n!} dy = e^{-c}.$$

Thus $P[zb_1 > b_1, \dots, zp_{n+1} > b_{n+1}] = e^{-b_1 - \dots - b_{n+1}}$, which means that zp_1, \dots, zp_{n+1} are independent and each has density function e^{-x} .

Noting that $\frac{zp_1}{zp_1 + \dots + zp_{n+1}} = p_1$, we have a result in the opposite direction:

Corollary. If y_1, \dots, y_{n+1} are independent, exponential random variables, then

$$\frac{y_1}{y_1 + \dots + y_{n+1}}, \frac{y_2}{y_1 + \dots + y_{n+1}}, \dots, \frac{y_{n+1}}{y_1 + \dots + y_{n+1}}$$

are uniformly distributed over the simplex

$$S_{n+1} = \{(x_1, \dots, x_{n+1}) : x_i \geq 0, \sum_{i=1}^{n+1} x_i = 1\}.$$

The representation of p_1, \dots, p_{n+1} in terms of independent exponential random variables y_1, \dots, y_{n+1} provides an elementary method for establishing well-known properties of uniform order statistics. For example, let

$$v_1 = \frac{y_1}{y_1 + \dots + y_{n+1}}, \quad v_2 = \frac{y_1 + y_2}{y_1 + \dots + y_{n+1}}, \quad \dots, \quad v_n = \frac{y_1 + \dots + y_n}{y_1 + \dots + y_{n+1}}.$$

Then (v_1, v_2, \dots, v_n) are uniformly distributed over the polytope $Q_n = \{(x_1, \dots, x_n) : 0 < x_1 < x_2 < \dots < x_n < 1\}$ and hence are distributed

as the order statistics of n independent uniform $[0,1]$ random variables.

Now let k be given, $1 \leq k \leq n$. Then v_1, \dots, v_{k-1} may be written

$$(5) \quad \frac{y_1}{y_1 + \dots + y_k} v_k, \quad \frac{y_1 + y_2}{y_1 + \dots + y_k} v_k, \quad \dots, \quad \frac{y_1 + \dots + y_{k-1}}{y_1 + \dots + y_k} v_k,$$

and v_{k+1}, \dots, v_n may be written

$$(6) \quad v_k + (1-v_k) \frac{y_{k+1}}{y_{k+1} + \dots + y_{n+1}}, \quad v_k + (1-v_k) \frac{y_{k+1} + y_{k+2}}{y_{k+1} + \dots + y_{n+1}}, \quad \dots, \\ v_k + (1-v_k) \frac{y_{k+1} + \dots + y_n}{y_{k+1} + \dots + y_{n+1}},$$

and it is obvious that, for given v_k , the random variables in (5) are distributed as the order statistics of $k-1$ uniform variables on the interval $[0, v_k]$, and are independent of the variables in (6), which are distributed as the order statistics of $(n-k)$ uniform variables on the interval $[v_k, 1]$.

One may prove directly that each of v_1, \dots, v_n has a beta distribution, but an interesting alternative method is to take advantage of known facts concerning the ratios of quadratic forms in normal variables, as each of the v 's may be viewed as such. The relationship between the beta and the F distributions may be brought in at this point, too.